



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2015

Spatio-phylogenetic multispecies distribution models

Kaldhusdal, Arne ; Brandl, Roland ; Müller, Jörg ; Möst, Lisa ; Hothorn, Torsten

Abstract: 1. Ecologists increasingly consider phylogenetic relatedness in both community composition and spatial arrangements in communities. Here we considered both the phylogenetic correlation between multiple species and the spatial correlation induced by unobserved spatial heterogeneity on multiple plots. For this analysis, we introduced phylogenetic spatial generalised linear mixed models (PSGLMMs), which are an extension of phylogenetic generalised linear mixed models (PGLMMs). 2. We used the framework of generalised linear array models to simultaneously model species and plot dimension. Such models have the potential to explain the correlation of the phylogenetic relationship of the observed species and of the spatial proximity of the plots, or both. We proposed model selection strategies based on proper scores and empirically evaluated them in a case study using bird count data. In our analysis, we focused on two special cases: the community composition model and the environmental sensitivity model. 3. Our simulation study indicated that it might be difficult to correctly identify phylogenetic signals when the phylogenetic correlation is rather low and when studying presence-absence or count data of rare or pervasive species.

DOI: <https://doi.org/10.1111/2041-210X.12318>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-115524>

Journal Article

Accepted Version

Originally published at:

Kaldhusdal, Arne; Brandl, Roland; Müller, Jörg; Möst, Lisa; Hothorn, Torsten (2015). Spatio-phylogenetic multispecies distribution models. *Methods in Ecology and Evolution*, 6(2):187-197.

DOI: <https://doi.org/10.1111/2041-210X.12318>

Spatio-Phylogenetic Multi-Species Distribution Models

Arne Kaldhusdal¹, Roland Brandl², Jörg Müller^{3,4}, Lisa Möst¹
and Torsten Hothorn⁵

¹Institut für Statistik, LMU München, Germany

²Fachbereich Biologie, Philipps-Universität Marburg, Germany

³Sachgebiet Forschung und Dokumentation,
Nationalparkverwaltung Bayerischer Wald, Grafenau, Germany

⁴Lehrstuhl für Terrestrische Ökologie, Technische Universität
München, Germany

⁵Institut für Epidemiologie, Biostatistik und Prävention,
Universität Zürich, Switzerland

Number of words abstract: 174

Number of words text: 6639

Number of figures: 3

Number of tables: 5

Number of references: 31

Number of online appendices: 2

Abstract

Ecologists increasingly consider phylogenetic relatedness in both community composition and spatial arrangements in communities. Here we considered both the phylogenetic correlation between multiple species and the spatial correlation induced by unobserved spatial heterogeneity on multiple plots. For this analysis, we introduced phylogenetic spatial generalised linear mixed models (PSGLMMs), which are an extension of phylogenetic generalised linear mixed models (PGLMMs). We used the framework of generalised linear array models to simultaneously model species and plot dimension. Such models have the potential to explain the correlation of the phylogenetic relationship of the observed species and of the spatial proximity of the plots, or both. We proposed model selection strategies based on proper scores and empirically evaluated them in a case study using bird count data. In our analysis, we focused on two special cases: the community composition model and the environmental sensitivity model. Our simulation study indicated that it might be difficult to correctly identify phylogenetic signals when the phylogenetic correlation is rather low and when studying presence-absence or count data of rare or pervasive species.

Keywords environmental gradient; generalised linear array models, null-model; phylogenetic community structure; phylogenetic signal; spatial corre-

23 lation

24 **1 Introduction**

25 With the increasing availability of phylogenetic information, two lines of
26 ecological and biogeographical inquiries have arisen: those leading to an un-
27 derstanding of adaptation (e.g. Hansen & Orzack, 2005) and those leading to
28 an understanding of processes that influence community composition and the
29 establishment of regional biotas (e.g. Lovette & Hochachka, 2006; Stephens
30 & Wiens, 2009). If we consider the evolution of species, e.g. as a Brownian
31 motion process (Felsenstein, 1985) in which genetic changes along an evolu-
32 tionary trajectory are random and generally small, the degree of divergence
33 between a pair of species should be proportional to the time since they di-
34 verged through speciation. In such cases, a phylogenetic signal is said to be
35 present between the species for the trait under study (Blomberg & Garland,
36 2002; Losos, 2008; Revell et al., 2008); for a clear description of the difference
37 between phylogenetic signal and phylogenetic niche conservatism see Losos
38 (2008). As traits mediate ecological processes, it follows that phylogenetic
39 relatedness within and between communities would allow these ecological
40 processes to be inferred. For this task, metrics have been developed to ex-
41 plore whether phylogenetically related species co-occur more often (e.g. by

42 environmental filtering) or less often (e.g. through competitive interactions)
43 than expected by chance (e.g. Webb et al., 2002; Vamosi et al., 2009). A
44 growing number of studies using metrics, such as the mean phylogenetic
45 distance between species or individuals, in combination with certain null-
46 models have shown that ecological communities are indeed phylogenetically
47 structured (reviewed, e.g. in Vamosi et al., 2009).

48 However, metrics like the mean phylogenetic distance between species or
49 individuals have a number of drawbacks (Ives & Helmus, 2011), such as
50 uncertainties of selecting the most appropriate and effective metric and null-
51 model for testing the hypothesis under consideration, low statistical power,
52 and the lack of possibilities for predictions. Therefore, Ives & Helmus (2011)
53 proposed model-based statistics in the form of phylogenetic generalised linear
54 mixed models (PGLMMs) to analyse the assembly of communities along
55 environmental gradients. The flexibility of the model allows one to design
56 tests for phylogenetic signals in occurrence data, and these tests have a higher
57 statistical power than tests using metrics and null-models. Furthermore,
58 these models provide possibilities to test for a phylogenetic signal in the
59 sensitivity to an environmental factor. Finally, PGLMMs can even consider
60 the case in which different species show similar reactions to environmental
61 factors, where, after consideration of this attraction, related species tend not
62 to co-occur (Ives & Helmus, 2011).

63 In our study, we embedded generalised linear mixed models for multiple
 64 species observed on multiple plots into the generalised linear array model
 65 (GLAM) framework (Currie et al., 2006). We show that PGLMMs as sug-
 66 gested by Ives & Helmus (2011) can be understood as a special case of
 67 GLAMs. Within the class of GLAMs, we extended PGLMMs to models that
 68 take into account both the phylogenetic correlation between species and the
 69 spatial correlation between plots. Within this new class of phylogenetic and
 70 spatial generalised linear mixed models (PSGLMMs), the correlation between
 71 the observations might be due to the phylogenetic correlation of the species,
 72 the spatial correlation of the plots, or both. Thus, efficient model-selection
 73 procedures are needed to differentiate between these three situations in a
 74 way driven by data, although careful interpretation is mandatory because
 75 the best fitting model might not describe the underlying process best. We
 76 suggest the application of cross-validated proper scores, such as the Dawid-
 77 Sebastiani or the Brier score, for model selection. On a more technical level,
 78 we applied standard software for generalised linear mixed model estimation
 79 (the R package `lme4`, Bates et al., 2013; R Core Team, 2013) for investigat-
 80 ing the empirical performance of the model selection procedures. Special
 81 emphasis was given to the development of tests for a phylogenetic signal in
 82 the co-occurrence pattern of species (Model I in Ives & Helmus, 2011) as
 83 well as tests for a phylogenetic signal in the sensitivity of species to envi-

ronmental factors (Model II in Ives & Helmus, 2011) for spatially correlated observations.

2 Methods

Our phylogenetic and spatial generalised linear mixed models (PSGLMMs) are GLMMs for which the random effects correlation structure is considered known and based on the phylogeny of the species considered. As a result, in addition to the fixed effects parameters, only a single variance parameter had to be estimated for each random effect modelled. The random effects variance parameter estimates were then used to decide whether a phylogenetic signal was present in the data. The model accepts any correlation structure derived from the phylogeny ranging from a correlation matrix based on the Brownian motion model with one free parameter to more complex models with more than one free parameter, like the accelerated (decelerated) model suggested by Blomberg et al. (2003).

The models may in principle be used for data with responses pertaining to any distribution belonging to the exponential family, but we focused on Gaussian, binary, and Poisson responses in our empirical investigations. For species (spp) $s = 1, \dots, S$ and plots $p = 1, \dots, P$, the responses were organised in an $(S \times P)$ -matrix \mathbf{Y} , i.e. in a two-dimensional array. Each column thus

103 represented the responses in a unique plot and each row represented the
 104 responses of a unique species. In the case of binary data, a plot-species
 105 datum y_{sp} took the value 1 if species s was present in plot p and took the
 106 value 0 if the species was absent. For count data, y_{sp} represented the number
 107 of observations made of species s in plot p , and for Gaussian data, any real
 108 value attribute/observation for species s in plot p may be used. Repeated
 109 measurements are possible, in which case \mathbf{Y} would be an array of dimension
 110 $(S \times P \times n)$, with n being the number of measurements per species and plot
 111 (incomplete designs can be filled with weight zero observations). For the sake
 112 of simplicity, we refer to the responses as occurrence data.

113 We assumed two different sources of variance in the species occurrence pat-
 114 terns. First, a species is more likely to occur in a plot if it is also present in
 115 neighbouring plots. This influence was determined by a spatial correlation
 116 matrix Σ_{plot} , based on some sort of neighbourhood/correlation structure of
 117 the plots. Second, a species is more likely to occur in a plot if a closely related
 118 species is also present (Ives & Helmus, 2011). This influence was determined
 119 by a phylogenetic correlation matrix Σ_{spp} . The combination of these two
 120 sources of variance will be referred to as spatio-phylogenetic variation.

121 Two different models were set up. The first model, the community com-
 122 position model (CCM), was used to test for a phylogenetic signal in the
 123 co-occurrence of species, i.e. whether their phylogenetic relationship has an

influence on the composition of communities. The second model, the environmental sensitivity model (ESM), was used to test for a phylogenetic signal in the sensitivity of a species to an environmental factor, i.e. whether phylogenetically closely related species respond to an environmental factor in a similar way. The CCM is thus the spatially enhanced version of Ives & Helmus' (2011) model I, and the ESM corresponds to their model II, but takes the spatial correlation of the plots into account.

Following Ives & Helmus (2011), we modelled phylogenetic attraction, where closely related species are more likely to co-occur, using a correlation matrix derived from a phylogenetic tree. To model phylogenetic repulsion, where closely related species are less likely to co-occur, we used the inverse of that correlation matrix as proposed by Ives & Helmus (2011).

2.1 Community Composition Model

In order to construct the design matrices for the CCM, we defined the species and plot matrices

$$\mathbf{S} = \mathbf{I}_S \in \mathbb{R}^{S \times S} \quad \text{and} \quad \mathbf{P} = \mathbf{I}_P \in \mathbb{R}^{P \times P},$$

where \mathbf{I}_S is an identity matrix of size S , and \mathbf{I}_P is an identity matrix of size P .

The corresponding phylogenetic and spatial correlation matrices, assumed to

141 be known and fixed, are denoted

$$\Sigma_{\text{spp}} \in \mathbb{R}^{S \times S} \quad \text{and} \quad \Sigma_{\text{plot}} \in \mathbb{R}^{P \times P}.$$

142 In order to account for the overall occurrence of the species, the CCM con-
 143 tained species-specific fixed effects. We accounted for the spatio-phylogenetic
 144 variation in the species occurrence by including random intercepts for each
 145 species and plot, which were assumed to be correlated in accordance with
 146 the phylogenetic and spatial correlation matrices. In the form of an array
 147 model (Currie et al., 2006), the CCM in its simplest form is

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{1}_P^\top + \mathbf{S}\mathbf{\Gamma}\mathbf{P}],$$

148 where $h[\cdot]$ is a response (i.e. inverse link) function suitable for the distribution
 149 of \mathbf{Y} (e.g., inverse logit, exponential, or identity) applied to each element
 150 of its argument. The matrix $\mathbf{B} \in \mathbb{R}^{S \times 1}$ contains the species fixed effects
 151 coefficients (as we only modelled a fixed intercept, \mathbf{B} is effectively a vector
 152 $\boldsymbol{\beta}_{\text{spp}} \in \mathbb{R}^S$ for the CCM). Because these are the same in all plots, the species
 153 matrix \mathbf{S} is implicitly expanded by the vector of P ones, $\mathbf{1}_P$, to form a
 154 design matrix. In this simple special case, the linear predictor can be written
 155 in the simpler form $\boldsymbol{\beta}_{\text{spp}}\mathbf{1}_P^\top + \mathbf{\Gamma}$. The matrix $\mathbf{\Gamma} \in \mathbb{R}^{S \times P}$ contains the random

156 effects coefficients for each species-plot datum. These were assumed to be
 157 distributed as

$$\text{vec}(\mathbf{\Gamma}) = \boldsymbol{\gamma} \sim \mathcal{N}_{SP}(\mathbf{0}, \sigma_{\gamma}^2 \boldsymbol{\Sigma}_{\text{plot}} \otimes \boldsymbol{\Sigma}_{\text{spp}}),$$

158 where the “vec” operator performs a column-wise concatenation of $\mathbf{\Gamma}$ and
 159 \otimes denotes the Kronecker product of two matrices. The spatio-phylogenetic
 160 variance parameter $\sigma_{\gamma}^2 = \sigma_{\text{plot}}^2 \sigma_{\text{spp}}^2$ is the only unknown random effects vari-
 161 ance parameter, and an estimate $\hat{\sigma}_{\gamma}^2 > 0$ would indicate the presence of a
 162 spatio-phylogenetic signal.

163 In order for standard mixed model algorithms to be applicable to the esti-
 164 mation of PSGLMMs, we transformed the random effects so that they were
 165 uncorrelated. We thus standardised them to be distributed as

$$\text{vec}(\mathbf{\Gamma}) = \boldsymbol{\gamma} \sim \mathcal{N}_{SP}(\mathbf{0}, \sigma_{\gamma}^2 \mathbf{I}_{SP})$$

166 while at the same time penalising the random effects design matrices ac-
 167 cordingly (this is a standard procedure in linear mixed models, see Fahrmeir
 168 et al., 2013). For this, we first required the Cholesky decomposition of the
 169 correlation matrix of the random effects. $\boldsymbol{\Sigma}_{\text{plot}}$ and $\boldsymbol{\Sigma}_{\text{spp}}$ were decomposed
 170 separately as

$$\mathbf{R}_{\Sigma_{\text{plot}}} \mathbf{R}_{\Sigma_{\text{plot}}}^{\top} = \Sigma_{\text{plot}} \quad \text{and} \quad \mathbf{R}_{\Sigma_{\text{spp}}} \mathbf{R}_{\Sigma_{\text{spp}}}^{\top} = \Sigma_{\text{spp}}.$$

171 As the Cholesky decomposition of a Kronecker product equals the Kronecker
 172 product of the Cholesky decomposition of each factor, we obtained the overall
 173 decomposition

$$(\mathbf{R}_{\Sigma_{\text{plot}}} \otimes \mathbf{R}_{\Sigma_{\text{spp}}})(\mathbf{R}_{\Sigma_{\text{plot}}} \otimes \mathbf{R}_{\Sigma_{\text{spp}}})^{\top} = \Sigma_{\text{plot}} \otimes \Sigma_{\text{spp}}.$$

174 Penalising the design matrix of the random effects by the factor $(\mathbf{R}_{\Sigma_{\text{plot}}} \otimes$
 175 $\mathbf{R}_{\Sigma_{\text{spp}}})$ allowed us to re-express the CCM as

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{1}_P^{\top} + \tilde{\mathbf{S}}\mathbf{\Gamma}\tilde{\mathbf{P}}], \tag{1}$$

176 using the penalised plot and species matrices

$$\tilde{\mathbf{P}} := \mathbf{P}\mathbf{R}_{\Sigma_{\text{plot}}} \quad \text{and} \quad \tilde{\mathbf{S}} := \mathbf{S}\mathbf{R}_{\Sigma_{\text{spp}}}.$$

177 Hereby, the species-plot specific intercepts $\mathbf{\Gamma}$ were implicitly standardised by
 178 the Cholesky factor $(\mathbf{R}_{\Sigma_{\text{plot}}} \otimes \mathbf{R}_{\Sigma_{\text{spp}}})^{-1}$.

179 2.2 Environmental Sensitivity Model

180 For the ESM, an environmental factor $\mathbf{x} \in \mathbb{R}^P$ was measured in each plot.
 181 We added fixed effects $\boldsymbol{\beta}_x \in \mathbb{R}^S$ for each species in addition to the fixed effects
 182 $\boldsymbol{\beta}_{\text{spp}}$, and thus assumed that each species also has a base sensitivity to the
 183 environmental factor. Because of the statistical difficulty in distinguishing
 184 between different phylogenetic effects (Ives & Helmus, 2011), the random
 185 intercepts used in the CCM were dropped in favour of random slopes in the
 186 environmental factor. This accounted for the spatio-phylogenetic variance
 187 associated with the sensitivity of the species to the environmental factor.
 188 We thus assumed a linear dependency between the environmental factor and
 189 the linear predictor of the ESM. When the plot matrix \mathbf{P} is replaced with an
 190 environmental factor matrix $\mathbf{X} = \text{diag}(\mathbf{x}) \in \mathbb{R}^{P \times P}$, a diagonal matrix with
 191 the values of \mathbf{x} on its diagonal, the ESM is expressed as

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{D}^\top + \tilde{\mathbf{S}}\mathbf{\Gamma}\tilde{\mathbf{X}}] = h[\boldsymbol{\beta}_{\text{spp}}\mathbf{1}_P^\top + \boldsymbol{\beta}_x\mathbf{x} + \tilde{\mathbf{S}}\mathbf{\Gamma}\tilde{\mathbf{X}}]. \quad (2)$$

192 As the ESM contains two fixed effects, $\mathbf{B} = \begin{bmatrix} \boldsymbol{\beta}_{\text{spp}} & \boldsymbol{\beta}_x \end{bmatrix} \in \mathbb{R}^{S \times 2}$ is now a
 193 matrix and \mathbf{S} is accordingly expanded by the matrix $\mathbf{D} = \begin{bmatrix} \mathbf{1}_P & \mathbf{x} \end{bmatrix} \in \mathbb{R}^{P \times 2}$
 194 to implicitly form a design matrix. We also defined the penalised environ-
 195 mental factor matrix as $\tilde{\mathbf{X}} := \mathbf{X}\mathbf{R}_{\Sigma_{\text{plot}}}$. $\mathbf{\Gamma}$ was assumed to be distributed as
 196 for the CCM. It is possible to model the sensitivity to several environmental

197 factors in one model by simply adding one random effect component for each
 198 factor. However, this would render the model selection step (Section 2.4)
 199 considerably more complex.

200 **2.3 Computational Considerations**

201 For the normal case, parameter estimation in linear mixed models is based
 202 on either the maximum likelihood estimator or the restricted maximum like-
 203 lihood estimator. In a first step, the (restricted) maximum likelihood esti-
 204 mator for σ_γ is computed by maximising the (restricted or penalised) profile
 205 log-likelihood. Formulating a PSGLMM as an array model helps to speed-up
 206 the evaluation of the linear predictor considerably, however, the optimisation
 207 procedure and thus also the estimates are not affected, and one can expect to
 208 obtain the exact same estimates through optimisation with or without using
 209 the array formulation. The estimation of the variance parameters σ_{plot}^2 and
 210 σ_{spp}^2 in models where the number of random effects equals the number of ob-
 211 servations (as in the PSGLMMs studied here) is especially challenging since
 212 σ_{plot}^2 and σ_{spp}^2 and the residual variance σ_ϵ are not well identifiable. In this
 213 situation, additional challenges arise and the traditional optimisers cannot
 214 be expected to converge well.

215 In the generalised case, where the conditional density of \mathbf{Y} is assumed to

216 follow a distribution from the one-parameter exponential family, model in-
 217 ference is based on the marginal density obtained by integration with respect
 218 to the random effects distribution. This integral is approximated and the
 219 quality of the estimation of σ_γ heavily depends on how close the marginal
 220 log-likelihood can be resembled. This process is completely independent of
 221 the array formulation and, as in the normal case, there will be considerable
 222 speed-ups to the evaluation of the approximate marginal log-likelihood but
 223 the estimates will be the same. There is no residual variance term in the
 224 logistic and Poisson models, so the variance σ_γ^2 is identifiable.
 225 In the context of anisotropic smoothing, Rodríguez-Álvarez et al. (2014) re-
 226 cently proposed an algorithm for the estimation of the variance parameters
 227 based on GLAMs. Although being much faster than a generalised mixed
 228 model, the results were virtually the same. So, overall improvements in con-
 229 vergence and accuracy can not be expected from the application of GLAMs
 230 alone.

231 **2.4 Model selection**

232 Even if the spatio-phylogenetic variance parameter σ_γ^2 is estimated to be
 233 non-zero, we have no guarantee that another correlation structure – even
 234 one not accounting for phylogeny – might not be better suited to describe

the variance in the observed data. The true correlation structure of the observations is not known, but four possible interesting scenarios are realistic for our setting: (1) the correlation is only phylogenetic, (2) the correlation is only spatial, (3) the correlation is spatio-phylogenetic, or (4) for Gaussian and Poisson responses, it is possible that the species-plot-specific variance is independent of both space and phylogeny. For binary responses, independent variances are not possible because of the constrained variance structure of binary data. For Gaussian responses in the CCM, the case of independent variance may also be neglected as it would be identical to the estimated residual variance. Selecting among different phylogenetic correlation structures allows an assessment of the Brownian motion assumption.

When testing for a phylogenetic signal in species community structure or in the sensitivity of the species to environmental factors, it is important to account for all reasonable correlation structures in the test. To test which correlation structure is actually the most suitable, we added one random effects component for each of the scenarios (1) through (3), and for (4) where applicable, to the CCM and ESM and let the model selection procedure decide which one is best suited to describe the covariance observed in the data.

For the CCM, we thus effectively fitted the full model

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{1}_P^\top + \underset{\text{phylogenetic}}{\tilde{\mathbf{S}}\mathbf{\Gamma}\mathbf{P}}^{(1)} + \underset{\text{spatial}}{\mathbf{S}\mathbf{\Gamma}\tilde{\mathbf{P}}}^{(2)} + \underset{\text{spatio-phylogenetic}}{\tilde{\mathbf{S}}\mathbf{\Gamma}\tilde{\mathbf{P}}}^{(3)} + \underset{\text{independent}}{\mathbf{S}\mathbf{\Gamma}\mathbf{P}}^{(4)}] \quad (3)$$

255 and for the ESM we fitted

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{D}^\top + \underset{\text{phylogenetic}}{\tilde{\mathbf{S}}\mathbf{\Gamma}\mathbf{X}}^{(1)} + \underset{\text{spatial}}{\mathbf{S}\mathbf{\Gamma}\tilde{\mathbf{X}}}^{(2)} + \underset{\text{spatio-phylogenetic}}{\tilde{\mathbf{S}}\mathbf{\Gamma}\tilde{\mathbf{X}}}^{(3)} + \underset{\text{independent}}{\mathbf{S}\mathbf{\Gamma}\mathbf{X}}^{(4)}]. \quad (4)$$

256 For example, for the second random effects component, (2) spatial, we did not
 257 need to penalise \mathbf{S} because we implicitly assumed a phylogenetic correlation
 258 matrix $\Sigma_{\text{spp}} = \mathbf{I}_S$. For each random effect component, a variance estimate
 259 $\hat{\sigma}_\gamma^2$ is obtained. Ideally only one of these estimates is non-zero, and the cor-
 260 relation structure pertaining to its random effect component is thus *chosen*.
 261 When all estimates are zero, this indicates a complete lack of group-specific
 262 variance across species and plots in the data. Our aim was to identify the
 263 correlation structure most likely to be present in the data. After selecting
 264 a correlation structure, a second model was fitted, using only the random
 265 effect component of the chosen correlation structure.
 266 If several variance estimates were positive, we thus needed to choose among
 267 these by model selection. Classical criteria for model selection, such as the
 268 AIC or BIC, are not adequate for model selection in the mixed models setting
 269 when the focus is on the choice of random effects (Vaida & Blanchard, 2005;

270 Braun et al., 2014). Vaida & Blanchard (2005) proposed a conditional AIC
 271 criterion for comparing linear mixed models with different random effects
 272 structures, based on inference on the conditional likelihood. The concept was
 273 extended to also apply to GLMMs by both Yu & Yau (2012) and Saefken
 274 et al. (2014). Another approach, suggested by Braun et al. (2014), uses mean
 275 cross-validated proper scores (see Gneiting & Raftery, 2007) to choose the
 276 model with the best predictive abilities in the GLMM setting. This method,
 277 which we utilised for the PSGLMs, is suitable for choosing random as well
 278 as fixed effects.

279 Because the predictive distribution F_y of the data y is not analytically ac-
 280 cessible, its first two central moments were obtained through a predictive
 281 cross-validation. For each plot p , we subsequently left one species s out and
 282 estimated its expectation μ_y and variance σ_y^2 based on the remaining species
 283 in plot p . For each plot-species datum, we calculated a proper score. For
 284 Gaussian and Poisson data, we applied the Dawid-Sebastiani score

$$S_{\text{DS}}(F_y, y) = -\frac{1}{2} \left[\log \sigma_y^2 + \frac{(y - \mu_y)^2}{\sigma_y^2} \right]$$

285 and for binary data, we applied the Brier score

$$S_{\text{B}}(F_y, y) = -(\mu_y - y)^2.$$

286 For each model containing only one random effect component, the mean of
 287 all proper scores was calculated and compared with that of the other models
 288 before a final choice was made. For both scores, a larger score means better
 289 predictive abilities. We refer to Braun et al. (2014) for technical details.
 290 The PSGLMMs may also be used to choose among several phylogenetic trees
 291 to decide which is best suited for the observed species data. If we assume
 292 that a spatio-phylogenetic signal is present, e.g. the CCM

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{1}_P^\top + \tilde{\mathbf{S}}^{(A)}\mathbf{\Gamma}\tilde{\mathbf{P}} + \tilde{\mathbf{S}}^{(B)}\mathbf{\Gamma}\tilde{\mathbf{P}}]$$

293 performs this test for two different phylogenetic correlation matrices Σ_{spp}^A
 294 and Σ_{spp}^B , derived from phylogenetic trees A and B. $\tilde{\mathbf{S}}^{(A)}$ and $\tilde{\mathbf{S}}^{(B)}$ are the
 295 species matrices penalised with the Cholesky factor corresponding to each
 296 phylogenetic correlation matrix.

297 **3 Simulation Study**

298 We assessed the performance of the PSGLMMs in simulation studies. For
 299 both the CCM and ESM, the three settings (a) *phylogenetic*, (b) *spatial*
 300 and (3) *spatio-phylogenetic* were simulated for each response type (Gaussian,
 301 binary, and Poisson). For all cases, responses were simulated for $S = 4$

species (spp) in $P = 100$ plots aligned on a 10×10 grid. The phylogenetic correlation matrix was set to

$$\Sigma_{\text{spp}} = \begin{bmatrix} 1 & \rho & 0 & 0 \\ \rho & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

where $\rho \in [0.02, 0.98]$ is the correlation between species 1 and 2. The responses of all other pairs of species were uncorrelated. The spatial correlation matrix Σ_{plot} was calculated using the Matérn correlation function, as given by Rue & Held (2005):

$$\frac{1}{2^{\kappa-1}\Gamma(\kappa)} \left(\frac{d}{\phi}\right)^{\kappa} K_{\kappa} \left(\frac{d}{\phi}\right),$$

with κ and ϕ both set to 1, and where d is the Euclidean distance between the centroids of any two plots and K_{κ} is the Bessel function. All plots were unit squares. To reduce computation time, all correlations < 0.10 were set to zero. This implied that the nearest 20 neighbours were used with a spatial correlation ranging from 0.11 to 0.37, independent of the phylogenetic correlation ρ and constant for all simulation runs.

For the ESM, we assumed that all base responses of species to the envi-

315 ronmental factor were zero. The first terms of the linear predictor matrices
 316 $\boldsymbol{\eta}_{\text{CCM}} = \boldsymbol{S}\boldsymbol{B}\mathbf{1}_P^\top + \boldsymbol{S}\boldsymbol{\Gamma}\boldsymbol{P}$ and $\boldsymbol{\eta}_{\text{ESM}} = \boldsymbol{S}\boldsymbol{B}\mathbf{1}_P^\top + \boldsymbol{S}\boldsymbol{\Gamma}\boldsymbol{X}$ were thus identical to
 317 $\boldsymbol{B} = \boldsymbol{\beta}_{\text{spp}}$, as given in Table 3. \boldsymbol{S} and \boldsymbol{P} were identity matrices of size S and
 318 P respectively and \boldsymbol{X} was a 10×10 array containing the values in each plot of
 319 the environmental factor \boldsymbol{x} . These were drawn from a P -dimensional normal
 320 distribution with mean vector $\boldsymbol{\mu}_x = \begin{bmatrix} 1 & 1 + 1/99 & \dots & 2 \end{bmatrix}^\top$ and covariance
 321 matrix $\boldsymbol{\Sigma}_x = 0.25 \cdot \boldsymbol{\Sigma}_{\text{plot}}$. The resulting values \boldsymbol{x} ranged from 0.27 to 3.05
 322 and exhibited a slight spatial pattern, with lower values in the lower left part
 323 and higher values in the upper right part of the array. The random effects $\boldsymbol{\Gamma}$
 324 were simulated in accordance with Table 1, and the responses obtained were
 325 as follows:

- 326 • Gaussian: $\text{vec}(\boldsymbol{Y}) \sim \mathcal{N}_{SP}(\text{vec}(\boldsymbol{\eta}), \sigma_\epsilon^2 \boldsymbol{I}_{SP})$ with residual variance $\sigma_\epsilon^2 =$
 327 0.01
- 328 • Binary: $\boldsymbol{Y}_{sp} \sim \mathcal{B}(1, h(\boldsymbol{\eta}_{sp}))$, where $h(\boldsymbol{\eta}_{sp}) = \mathbb{P}[\boldsymbol{Y}_{sp} = 1]$ and the re-
 329 sponse function is $h(\boldsymbol{\eta}) = \exp\{\boldsymbol{\eta}\} / (1 + \exp\{\boldsymbol{\eta}\})$
- 330 • Poisson: $\boldsymbol{Y}_{sp} \sim \mathcal{P}(h(\boldsymbol{\eta}_{sp}))$, where $h(\boldsymbol{\eta}_{sp}) = \mathbb{E}[\boldsymbol{Y}_{sp}]$ and $h(\boldsymbol{\eta}) = \exp\{\boldsymbol{\eta}\}$

331 with $s = 1, \dots, S$ and $p = 1, \dots, P$.

332 [Table 1 about here.]

333 In the phylogenetic setting (a), the correlation of the simulated responses

334 thus increased with ρ as it depended on the phylogenetic correlation matrix
 335 Σ_{spp} . In the spatial setting (b), the correlation of the responses was constant
 336 for all ρ , as it only depended on the spatial correlation matrix Σ_{plot} . In
 337 the spatio-phylogenetic setting (c), the correlation rose with ρ , but was also
 338 dependent on the spatial correlation matrix.

339 The phylogenetic and spatio-phylogenetic random effect components (1) and
 340 (3) in equations 3 and 4 also varied with ρ as the penalised species matrix
 341 $\tilde{\mathbf{S}}$ depended on the Cholesky factor of Σ_{spp} , i.e. the same phylogenetic cor-
 342 relation ρ used for generating the data also defined Σ_{spp} in the model. The
 343 spatial and independent random effect components (2) and (4) were fixed
 344 for all ρ . For each simulation setting, the *correct* random effect component
 345 chosen by the model was thus the one carrying the same name as the setting.

346 For each response type and setting, 100 data sets were simulated. The re-
 347 sulting parameter estimates are found in Figure 1 and 3 and 4 in the online
 348 supplement for the CCMs and in Figures 5–7 in the online supplement for
 349 the ESMs. The results for each simulation setting are given in separate
 350 columns, with one plot for each random effect component modelled. For low
 351 correlations ρ , Σ_{spp} was similar to the identity matrix \mathbf{I}_G . In the phyloge-
 352 netic simulation setting (a), the phylogenetic and independent random effect
 353 components were thus also similar for low correlations, and we expected the
 354 estimated variance to be split between them. This was also the case for the

spatial and spatio-phylogenetic random effect components in the spatial and spatio-phylogenetic settings (b) and (c). For increasing correlations, we expected the variance estimate pertaining to the correct component, i.e. the one used to simulate the data in that setting, to level off around the true parameter, while all other variance estimates were expected to tend towards zero.

Not all simulated runs could be evaluated. Especially for Gaussian responses (between 16% and 19% of all runs were discarded for the CCM and between 25% and 29% for the ESM) and for Poisson responses for the ESM (between 2% and 5% of all runs were discarded), a considerable number were discarded (see Table 2). For Gaussian responses, all discarded runs were due to the parameter estimation (in `lme4`) not converging due to nonidentifiable variance and residual variance parameters, see Section 2.3. For Poisson responses for the ESM, most discarded runs were due to numerical issues pertaining to the predictive cross-validation, see Braun et al. (2014).

[Table 2 about here.]

3.1 Results CCM

The variance parameter estimates for the CCM for Gaussian responses in Figure 1 all clearly followed the expected trends. For low values of ρ in the

374 phylogenetic setting, the estimated phylogenetic variance $\hat{\sigma}_{\text{spp}}^2$ and the resid-
 375 ual variance $\hat{\sigma}_{\epsilon}$ were close to one, although the corresponding true parameters
 376 were $\sigma_{\text{spp}}^2 = 2$ and $\sigma_{\epsilon} = 0.01$ as a consequence of near nonidentifiability. In
 377 the spatial setting (b), the results depended on ρ only indirectly via the phy-
 378 logenetic term in the model. Low values of ρ rendered parameter estimation
 379 more difficult and were associated with a higher variability of $\hat{\sigma}_{\text{plot}}^2$. It is
 380 plausible that the estimates in the spatio-phylogenetic setting (c) were less
 381 precise than the estimates in the phylogenetic (a) and spatial (b) settings,
 382 as the variance parameter estimated was made up of two separate terms
 383 ($\sigma_{\text{spp}}^2 \sigma_{\text{plot}}^2 = 2 \cdot 2 = 4$). For all cases, the parameter estimates levelled off at
 384 values quite close to the true parameters (see Table 3).

385 [Table 3 about here.]

386 [Figure 1 about here.]

387 For binary responses, the expected trends were less apparent (see Figure 3 in
 388 online supplement). A positive random effect variance parameter estimate in
 389 the phylogenetic setting (a) was not notable until the correlation ρ reached
 390 a medium strength, but the estimate rose steadily. This positive trend was
 391 also clear for the spatio-phylogenetic setting (c); in the spatial setting (b),
 392 the estimations were split between the spatial and spatio-phylogenetic com-
 393 ponents for all correlations. This was caused by the relatively weak spatial

394 correlation structure Σ_{plot} used. All variance estimates were highly variable.
 395 The results for Poisson responses (see Figure 4 in online supplement) were
 396 more or less as in the case of Gaussian responses, except with considerably
 397 more variability to the parameter estimates. The parameter estimates also
 398 tended more slowly towards the true parameters than estimates for Gaussian
 399 responses.

400 Figure 2 depicts the selection rates of the correct random effect components,
 401 based on the predictive cross-validation, for each response type and simula-
 402 tion setting for the CCM. The red line illustrates the selection rate of the
 403 correct random effect component; for example, in the phylogenetic setting
 404 (a), it indicates the selection rate of the random effect component phyloge-
 405 netic. The blue line illustrates the rate with which any of the two components
 406 in the phylogenetic or spatio-phylogenetic setting were selected, i.e. at which
 407 rate any phylogenetic signal was detected. In the spatial setting (b), the
 408 latter rate is thus false positive and should ideally be low.

409 [Figure 2 about here.]

410 All selection rates for Gaussian responses reached a very high level for rela-
 411 tively low correlations ρ . For the spatio-phylogenetic setting (c), however, a
 412 severe drop in the selection rate for high correlations contradicted the ten-
 413 dency of the variance parameter estimates in Figure 1. As can be seen from

414 the blue curve, this was due to the phylogenetic component being chosen
415 instead.

416 The selection rates for binary responses also reached a relatively high level
417 for the phylogenetic and spatio-phylogenetic settings (a) and (c). For the
418 spatial setting (b), however, the inability of the model to choose between
419 the spatial and spatio-phylogenetic components became even clearer, with a
420 selection rate constantly at around 0.50. The results indicate that it is in
421 fact rather difficult for the model and model selection procedure to decide
422 which correlation structure (spatio-phylogenetic or spatial) fits the data best,
423 when in fact only a moderately large spatial correlation is present in binary
424 or count data as was the case in our simulation study.

425 For Poisson responses as well, the selection rate of the correct component
426 in the spatial setting (b) was low. This trend was due to the phylogenetic
427 and spatio-phylogenetic components being chosen for Gaussian and binary
428 responses, whereas for Poisson responses, the independent component (4)
429 was chosen most often. In the spatio-phylogenetic setting (c), only a medium
430 selection rate of around 0.50 was reached for the correct component, although
431 the selection rate of either of the two components, phylogenetic and spatio-
432 phylogenetic, was relatively high for medium to strong correlations.

433 **3.2 Results ESM**

434 For the ESMs, the trends of the parameter estimates were more precise and
435 clearer than for the CCM for all response types, but otherwise did not differ
436 noticeably (see Figures 5–7 in the online supplement).

437 The ESM selection rates (Figure 3) differed from those of the CCM. While
438 the selection rate in the phylogenetic setting (a) was somewhat higher for
439 binary responses, it was considerably lower for Poisson responses for medium
440 to high correlations. On the other hand, in the spatio-phylogenetic setting
441 (c), the selection rate was much better for Poisson responses and somewhat
442 worse for Gaussian responses. The selection rate for Gaussian responses in
443 the spatial setting (b) was also much lower than for the CCM. However, the
444 false positive rate was approximately the same.

445 [Figure 3 about here.]

446 **3.3 Comparison to PGLMMs**

447 To compare the performance of the original implementation of the phyloge-
448 netic generalised linear mixed models (PGLMMs) of Ives & Helmus (2011)
449 with our `lme4`-based implementation of the PSGLMMs, we fitted both mod-
450 els to 100 simulated binary occurrence data sets and compared the results.
451 As the models in the special case of spatially uncorrelated data are the same,

452 differences can be attributed to different software implementations used. The
 453 estimation technique applied for the PGLMMs by Ives & Helmus' (2011) in
 454 the R package `picante` (Kembel et al., 2010) is a combination of penalised
 455 quasi-likelihood (PQL) and REML estimation. For the PSGLMMs, which
 456 make use of the R package `lme4` (Bates et al., 2013), the random effect vari-
 457 ance parameters are estimated via Laplace approximation.
 458 The data were simulated using the function `pglmm.sim()` from `picante`;
 459 the corresponding simulation model is described in Ives & Helmus (2011).
 460 For each simulation run, the occurrence of $S = 16$ species was simulated in
 461 $P = 30$ plots. During simulation, plots with fewer than two species present
 462 were discarded. The actual number of plots may thus be smaller and may dif-
 463 fer between the simulation runs. The phylogenetic relationship of the species
 464 was characterised through a balanced phylogenetic tree, from which a phy-
 465 logenetic correlation matrix Σ_{spp} was derived. The species-specific random
 466 effects were simulated using a variance estimate $\sigma_{\text{spp}}^2 = 1$. A random ef-
 467 fect accounting for differences in the number of species in each plot was also
 468 added, with a small but positive variance parameter σ_{plot}^2 .
 469 In array form, the model was set up as

$$\mathbb{E}[\mathbf{Y}] = h[\mathbf{S}\mathbf{B}\mathbf{1}_P^\top + \tilde{\mathbf{S}}\mathbf{\Gamma}\mathbf{P} + \mathbf{S}\mathbf{\Gamma}\mathbf{P}],$$

470 with random effects $\mathbf{\Gamma} \sim \mathcal{N}_{SP}(0, \sigma_\gamma^2 \mathbf{I}_P \otimes \mathbf{I}_S)$, fixed effects $\mathbf{B} = \boldsymbol{\beta}_{\text{spp}}$ and
 471 $\tilde{\mathbf{S}} = \mathbf{S} \mathbf{R}_{\Sigma_{\text{spp}}}$. \mathbf{S} and \mathbf{P} are identity matrices of size S and P . The first
 472 random effect component models the random effects of the species, while
 473 accounting for their phylogenetic relatedness. The second component models
 474 an unstructured random effect in the plots.

475 In their original model, Ives & Helmus (2011) assume a variance-covariance
 476 structure $\sigma_{\text{plot}}^2 \mathbf{I}_P \otimes \mathbf{J}_S$ for the random effect in the plots, where \mathbf{J}_S is an
 477 $(S \times S)$ -matrix of ones. Because this matrix is not positive semi-definite it
 478 cannot be used for our PSGLMMs. In order to allow for the results from
 479 the two models to be compared we instead utilised the variance-covariance
 480 structure $\sigma_{\text{plot}}^2 \mathbf{I}_P \otimes \mathbf{I}_S$ for both models.

481 The variance parameters σ_{spp}^2 and σ_{plot}^2 were estimated for the PGLMM using
 482 the function `pglmm.fit()` from `picante`. For the PSGLMM, we used the
 483 code given in the online supplements, which made use of `glmer()` from `lme4`.

484 In addition to obtaining the random effect variance estimates with both meth-
 485 ods, we also calculated the REML log-likelihood and considered the fixed
 486 effects obtained with the two different methods. As the estimation tech-
 487 nique was the only difference between the fitted models, we expected their
 488 results, presented in Figures 1 and 2 in the online supplement, to be similar.

489 Note that 28 runs were discarded, either because the PGLMM (5) or the
 490 PSGLMM (18) estimation algorithm failed to converge or because of issues

491 of quasi-complete separation (5).
 492 The log-likelihood was practically identical for both models in all considered
 493 runs. The variance estimates $\hat{\sigma}_{\text{spp}}^2$ were in general higher and varied more for
 494 the PSGLMM than for the PGLMM, but the estimates of the two models
 495 were highly correlated. The median estimate was 0.578 for the PGLMM
 496 and 0.788 for PSGLMM. Regarding the estimates $\hat{\sigma}_{\text{plot}}^2$ for the PGLMM and
 497 PSGLMM, we observed no coherence other than most estimates being close
 498 to zero. The fixed effect parameter estimates were also highly correlated
 499 between the two models, although they were generally somewhat lower for
 500 the PSGLMM.

501 **3.4 Summary of Simulation Results**

502 In general, PSGLMMs based on their implementation in the package `lme4`
 503 were able to estimate the fixed and variance parameters and the suggested
 504 model selection procedure performed well in selecting the correct model in
 505 spatio-phylogenetically correlated Gaussian, binary, and Poisson data.
 506 Exceptions were observed in the following situations: For Gaussian responses
 507 with a low phylogenetic correlation, the variance parameters pertaining to
 508 phylogenetic and spatio-phylogenetic model terms and the residual variance
 509 were not identifiable resulting in biased estimates with a large variability.

510 Although we did not alter the spatial correlation here, the same would apply
511 to low spatial correlations in a data set. Therefore, one cannot expect the
512 procedure to work well for low phylogenetic or spatial correlations. For higher
513 correlations ($\rho > 0.1$ in our simulations), the model performed very well in
514 all three scenarios and the model selection procedure led to the correct model
515 with high probability, although it seemed to be hard to differentiate between
516 phylogenetic and spatio-phylogenetic terms for very large correlations.

517 For binary responses, the variability of the estimated variance parameters
518 was much higher, due to the lower information contained in binary data. As
519 a consequence, the selection frequency of the correct model was not quite
520 as high as in the Gaussian case but still acceptable. Note that the low
521 selection frequency in the spatial setting (b) was due to the relatively low
522 spatial correlation used in our simulations. For more strongly correlated
523 observations, the corresponding selection frequency would be much higher.

524 The same applies to the Poisson case.

525 Our simulation results for the phylogenetic setting (a) carry over to the
526 traditional PGLMM model because the two software implementations (based
527 on the `nlme` and `lme4` packages) led to practically identical results.

528 4 Application to Bird Abundance Data

529 In the following section, the PSGLMs were applied to a data set of bird
530 counts to test for phylogenetic signals. The data consist of observational
531 counts of 47 bird species belonging to the orders *Piciformes* and *Passer-*
532 *iformes*, made along four transects through the Bavarian Forest National
533 Park in Germany (see Figure 8 in the online supplement). Counts ranging
534 from 0 to 37 were made on a total of 371 plots. With 74% of the data being
535 zeros, the empirical distribution was strongly positively skewed.

536 We chose a subgroup of six closely related species consisting of the passer-
537 ine (*Passeriformes*) species *Poecile montanus*, *Periparus ater*, *Lophophanes*
538 *cristatus*, *Poecile palustris*, *Parus major* and *Cyanistes caeruleus*, with a
539 mean phylogenetic correlation of 0.808 between the species. The spatial
540 correlation matrix was calculated on the basis of the Euclidean distances be-
541 tween the plots using the Matérn correlation function. To reduce computa-
542 tional time, spatial correlations < 0.05 were set to zero. For the ESM, we con-
543 sidered five environmental factors: elevation, coverage (in %) of the middle
544 forest layer by broadleaved trees (in the following referred to as *broadleaved*
545 *trees*) and by rejuvenating broadleaved trees (*rejuvenation*), the maximum
546 trunk diameter (in *cm*) at breast height (*DBH*) and the total number of tree
547 holes within an area of 0.1ha (*tree holes*).

548 4.1 Results

549 In addition to fitting a CCM, we fitted one ESM for each environmental
550 factor. To all models, we added the six random effect components listed
551 in Table 4. The component *spatial* covers the case where the variance is
552 of a purely spatial nature – independent of the phylogenetic relatedness of
553 the birds. To cover the cases where the variance is dependent on only phy-
554 logeny or on phylogeny and space, we added the components *phylogenetic*
555 and *spatio-phylogenetic* twice – once for an attraction tendency and once for
556 a repulsion tendency (see Ives & Helmus, 2011). For an attraction tendency,
557 we assumed that closely related species are more likely to co-occur than more
558 distantly related species. We utilised the phylogenetic correlation matrix to
559 incorporate this case. For a repulsion tendency, we assumed that closely re-
560 lated species are less likely to co-occur than more distantly related species.
561 Following Ives & Helmus (2011) we utilised the inverse of the phylogenetic
562 correlation matrix, Σ_{spp}^{-1} , for this case. Lastly, we also added the component
563 *independent* to cover the possibility of an observed variance in the sensitivity
564 of the species to the environmental factors to be independent of space and
565 phylogeny.

566 [Table 4 about here.]

567 The choice of a random effect component by the PSGLMM thus allowed us
 568 to decide whether or not a phylogenetic signal was detected in the sensitivity
 569 of the species to an environmental factor, and if so, whether closely related
 570 species have a tendency to attract or repel each other. Fixed intercepts were
 571 added for each species to account for their overall abundance. To account
 572 for the different base sensitivity of the species to the environmental factors,
 573 fixed slopes in the environmental factors were also added for each species.
 574 For the CCM, we expected either the phylogenetic or the spatio-phylogenetic
 575 component with an attraction tendency to be chosen, and this was the case:
 576 a spatio-phylogenetic community structure in which closely related species
 577 are more likely to co-occur (attraction) was chosen by the CCM with an
 578 estimated random effect variance of $\hat{\sigma}_{\gamma}^2 = 1.860$. Because of the spatial
 579 patterns of the land-cover types, we expected either the spatial or one of the
 580 spatio-phylogenetic random effect components to be chosen for the ESMs.
 581 We further expected an attraction tendency as these birds species are not
 582 known to compete for resources.

583 A spatio-phylogenetic signal was detected in the sensitivity of the species to
 584 tree holes, broadleaved trees, and rejuvenation (see Table 5). For the latter
 585 two environmental factors, the species exhibited an attraction tendency, i.e.,
 586 the closer the species were related, the more likely they were to share similar
 587 values of these environmental factors. For tree holes, the species exhibited

588 a repulsion tendency, with more distantly related species being more likely
589 to share similar values. For the environmental factors DBH and elevation, a
590 purely spatial correlation structure was chosen. The species were thus more
591 likely to share similar values when spatially close to each other, regardless of
592 their phylogenetic relatedness.

593 [Table 5 about here.]

594 **5 Conclusion**

595 Numerous studies have already provided convincing evidence that ecological
596 communities exhibit phylogenetic structure at both small and large spatial
597 scales (e.g. Graves & Gotelli, 1993; Webb et al., 2002; Lovette & Hochachka,
598 2006; Helmus et al., 2007; Pillar & Duarte, 2010; Riedinger et al., 2013).
599 Most of these authors used one or several of the available metrics (e.g. mean
600 phylogenetic distance between species or individuals) to summarise the phy-
601 logenetic structure as well as various null-models to explore whether the
602 selected metric deviates from the expectation of a random distribution of
603 species. But as noted by Ives & Helmus (2011), such metrics summarise the
604 complexities of the community structure in a single number and therefore
605 only allow for tests of basic hypotheses. Nevertheless, the use of multiple
606 metrics in a study may increase our confidence in the analyses and interpre-

607 tations as long as all metrics indicate the same pattern. In the case where
 608 the metrics signal different patterns, we have the problem either to dismiss
 609 the pattern or to make an often arbitrary selection of a metric we want to
 610 rely on for further interpretations. To overcome these limitations, Ives &
 611 Helmus (2011) have implemented generalised linear mixed models for phylo-
 612 genetic community analyses. They have shown that such mixed models have
 613 a higher statistical power in detecting phylogenetic signals than approaches
 614 that rely on metrics and null-models.

615 Considering PGLMMs as special cases of GLAMs has three merits: (1) the
 616 extension to more complex models, such as one allowing a spatial correlation
 617 between observations, is possible; (2) computationally attractive properties
 618 of GLAMs will carry over to P(S)GLMMs; and (3) novel model selection
 619 procedures can be applied. Although Currie et al. (2006) present a proto-
 620 type implementation of generalised linear mixed models, no production-ready
 621 code is currently available. We therefore used the `lme4` package to implement
 622 two of the models proposed by Ives & Helmus (2011): one model used the
 623 qualitative or quantitative composition of the community to test for a phylo-
 624 genetic structure in a community matrix (Model I of Ives & Helmus, 2011),
 625 and the other model tested for a phylogenetic signal in the sensitivity of the
 626 species to environmental factors. Although the approaches are not identical,
 627 a comparison of the model fits using our PSGLMM and the implementation

628 of Ives & Helmus' (2011) PGLMM using simulated data indicated that the
629 two approaches indeed perform on par in this special situation. Hence, our
630 empirical findings in the scenario phylogenetic (a) carry over to the corre-
631 sponding implementation of the phylogenetic models described by Ives &
632 Helmus (2011). However, the novel model formulation presented here allows
633 the spatial correlation structure between plots to be taken into account. In
634 addition, selection between various correlation structures was also introduced
635 for PSGLMMs.

636 PSGLMMs are very complex models relying on a number of rather strict
637 assumptions. All limitations of the generalised linear mixed model frame-
638 work apply, such as the conditional distribution of the response coming from
639 an exponential family, conditional independence of the observations, nui-
640 sance parameters being independent of explanatory variables, and indepen-
641 dent normal random effects. The latter assumption might be hard to justify
642 since neither the species- nor plot-specific random effects can be assumed to
643 be symmetric in general. The ability to let the model select one of several
644 phylogenetic correlation structures relaxes strict assumptions on the phylo-
645 genetic signal. In principle, this would also be possible for the spatial cor-
646 relation, however, the Matérn family is one of the most flexible approaches
647 and a standard choice in spatial statistics (Fahrmeir et al., 2013).

648 Our simulation studies showed that PSGLMMs are generally capable of se-

649 lecting the correct random effect components and yield precise estimates of
 650 the variance parameters. Although the selection rate of the correct ran-
 651 dom effect component was somewhat low for low correlations, the models
 652 were effective in detecting whether a phylogenetic signal (spatial or not) was
 653 present. However, both the CCM and the ESM exhibited a quite high false-
 654 positive rate for binary and Poisson responses. For binary occurrence data,
 655 it is important that rare or pervasive species are not included when fitting
 656 the models as these often lead to issues concerning quasi-complete separa-
 657 tion, which leads to estimates of very large positive or negative fixed effects.
 658 For Poisson responses, especially rare species are problematic. A technical
 659 consequence of the large negative fixed effects is that they render the pre-
 660 dictive cross-validation infeasible, as weights pertaining to these species are
 661 set to zero and later used as denominators. Overall, one cannot expect PS-
 662 GLMMs to identify the correct correlation structure (phylogenetic, spatial,
 663 or spatio-phylogenetic) for low correlations. When the correlations can be
 664 assumed to be rather large, PSGLMMs and the model selection procedure
 665 have the potential to select the correlation structure that comes closest to
 666 reality. Improved optimisers able to deal with the near-nonidentifiability of
 667 the variance parameters in the Gaussian case may offer improvements in the
 668 situation of low phylogenetic or spatial correlations.
 669 The computational restrictions in our study considering only six species can

670 be relaxed by utilising the array structure of the data. This would reduce
671 computational time drastically and would allow us to fit the models to large
672 community matrices in situations where convergence of the mixed model es-
673 timation is unproblematic. Because of the small number of species used, our
674 example should be considered as an illustration. Nevertheless, the applica-
675 tion to the bird count data in Bavaria revealed spatio-phylogenetic signals
676 both in the community structure and in the species' sensitivity to several
677 environmental factors. Thus, the application of PSGLMMs allowed the de-
678 tection of sophisticated spatial-phylogenetic patterns.

679 **Acknowledgements**

680 We thank Karen A. Brune for improving the language.

681 **Data Accessibility**

682 Data and computer code for data analysis and simulation experiments are
683 contained in the R add-on package **TH.data**, version 1.0-5 or higher, avail-
684 able from <http://CRAN.R-project.org/package=TH.data>.

685 References

- 686 Bates, D.; Maechler, M.; Bolker, B. & Walker, S. (2013) *lme4: Linear mixed-*
687 *effects models using Eigen and S4*, R package version 1.0-0.
- 688 Blomberg, S.P. & Garland, T. (2002) Tempo and mode in evolution: phy-
689 logenetic inertia, adaptation and comparative methods, *Journal of Evolu-*
690 *tionary Biology*, 15(6), 899–910.
- 691 Blomberg, S.P.; Garland, T. & Ives, A.R. (2003) Testing for phylogenetic
692 signal in comparative data: Behavioral traits are more labile, *Evolution*,
693 57(4), 717–745.
- 694 Braun, J.; Sabanés Bové, D. & Held, L. (2014) Choice of generalized linear
695 mixed models using predictive crossvalidation, *Computational Statistics &*
696 *Data Analysis*, 75, 190–202.
- 697 Currie, I.D.; Durban, M. & Eilers, P.H.C. (2006) Generalized linear array
698 models with applications to multidimensional smoothing, *Journal of the*
699 *Royal Statistical Society Series B*, 68(2), 259–280.
- 700 Fahrmeir, L.; Kneib, T.; Lang, S. & Marx, B. (2013) *Regression: Models,*
701 *Methods and Applications*, Springer.

- 702 Felsenstein, J. (1985) Phylogenies and the comparative method, *The Amer-*
703 *ican Naturalist*, 125(1), 1–15.
- 704 Gneiting, T. & Raftery, A.E. (2007) Strictly proper scoring rules, prediction,
705 and estimation, *Journal of the American Statistical Association*, 102, 359–
706 378.
- 707 Graves, G.R. & Gotelli, N.J. (1993) Assembly of avian mixed-species flocks in
708 Amazonia, *Proceedings of the National Academy of Sciences of the United*
709 *States of America*, 90(4), 1388–1391.
- 710 Hansen, T.F. & Orzack, S.H. (2005) Assessing current adaptation and phy-
711 logenetic inertia as explanations of trait evolution: The need for controlled
712 comparisons, *Evolution*, 59(10), 2063–2072.
- 713 Helmus, M.R.; Savage, K.; Diebel, M.W.; Maxted, J.T. & Ives, A.R. (2007)
714 Separating the determinants of phylogenetic community structure, *Ecology*
715 *Letters*, 10(10), 917–925.
- 716 Ives, A.R. & Helmus, M.R. (2011) Generalized linear mixed models for phy-
717 logenetic analyses of community structure, *Ecological Monographs*, 81(3),
718 511–525.
- 719 Kembel, S.; Cowan, P.; Helmus, M.; Cornwell, W.; Morlon, H.; Ackerly, D.;

- 720 Blomberg, S. & Webb, C. (2010) Picante: R tools for integrating phyloge-
721 nies and ecology, *Bioinformatics*, 26, 1463–1464.
- 722 Losos, J.B. (2008) Phylogenetic niche conservatism, phylogenetic signal and
723 the relationship between phylogenetic relatedness and ecological similarity
724 among species, *Ecology Letters*, 11(10), 995–1003.
- 725 Lovette, I.J. & Hochachka, W.M. (2006) Simultaneous effects of phyloge-
726 netic niche conservatism and competition on avian community structure,
727 *Ecology*, 87(7), 14–28.
- 728 Pillar, V.D. & Duarte, L.d.S. (2010) A framework for metacommunity anal-
729 ysis of phylogenetic structure, *Ecology Letters*, 13(5), 587–596.
- 730 R Core Team (2013) *R: A language and environment for statistical computing*,
731 R Foundation for Statistical Computing, Vienna, Austria.
- 732 Revell, L.J.; Harmon, L.J. & Collar, D.C. (2008) Phylogenetic signal, evolu-
733 tionary process, and rate, *Systematic Biology*, 57(4), 591–601.
- 734 Riedinger, V.; Müller, J.; Stadler, J.; Ulrich, W. & Brandl, R. (2013) Assem-
735 blages of bats are phylogenetically clustered on a regional scale, *Basic and*
736 *Applied Ecology*, 14(1), 74–80.
- 737 Rodríguez-Álvarez, M.X.; Lee, D.J.; Kneib, T.; Durbán, M. & Eilers, P.

- 738 (2014) Fast smoothing parameter separation in multidimensional general-
739 ized P-splines: the SAP algorithm, *Statistics and Computing*, in press.
- 740 Rue, H. & Held, L. (2005) *Gaussian Markov random fields: Theory and*
741 *applications*, vol. 104 of *Monographs on Statistics and Applied Probability*,
742 Chapman & Hall, London.
- 743 Saefken, B.; Kneib, T.; van Waveren, C.S. & Greven, S. (2014) A unifying
744 approach to the estimation of the conditional Akaike information in gen-
745 eralized linear mixed models, *Electronic Journal of Statistics*, 8, 201–225.
- 746 Stephens, P.R. & Wiens, J.J. (2009) Bridging the gap between community
747 ecology and historical biogeography: niche conservatism and community
748 structure in emydid turtles, *Molecular Ecology*, 18(22), 4664–4679.
- 749 Vaida, F. & Blanchard, S. (2005) Conditional Akaike information for mixed
750 effects models, *Biometrika*, 92(2), 351–370.
- 751 Vamosi, S.M.; Heard, S.B.; Vamosi, J.C. & Webb, C.O. (2009) Emerging
752 patterns in the comparative analysis of phylogenetic community structure,
753 *Molecular Ecology*, 18(4), 572–592.
- 754 Webb, C.O.; Ackerly, D.D.; McPeck, M.A. & Donoghue, M.J. (2002) Phylo-
755 genies and community ecology, *Annual Review of Ecology and Systematics*,
756 33, 475–505.

757 Yu, D. & Yau, K.K.W. (2012) Conditional Akaike information criterion for
758 generalized linear mixed models., *Computational Statistics & Data Analy-*
759 *sis*, 56(3), 629–644.

760 List of Tables

761	1	Distributions used to simulate the random effects for the dif-	
762		ferent correlation structures in each setting (a) through (c). . .	46
763	2	The number of discarded runs in the simulation study (from	
764		a total of 4,900 runs per response type, setting and model). . .	47
765	3	Fixed effects parameters β_{spp} for each of the four species and	
766		phylogenetic and spatial random effect variance parameters,	
767		σ_{spp}^2 and σ_{spp}^2 respectively, used for each response type in the	
768		simulation study of the PSGLMs.	48
769	4	Variance-covariance matrices for each random effect compo-	
770		nent used to model the bird abundance data. For a <i>repulsion</i>	
771		tendency, where closely related species are less likely to co-	
772		occur, contrary to the attraction tendency, the inverse of the	
773		phylogenetic correlation matrix was used.	49
774	5	For each of five environmental factors, a separate ESM was	
775		fitted. With the spatial correlation structure chosen for DBH	
776		and elevation, we concluded that no phylogenetic signal is	
777		present in the response of the species to the factors. For the	
778		other factors, a spatio-phylogenetic signal was detected with a	
779		repulsion tendency for tree holes and an attraction tendency	
780		for broadleaved trees and rejuvenation. $\hat{\sigma}_{\gamma}^2$ is the random effect	
781		variance estimate for the correlation structure chosen using a	
782		predictive cross-validation.	50

CORRELATION STRUCTURE	$\text{vec}(\mathbf{\Gamma}) \sim$	SIMULATION
(1) phylogenetic	$\mathcal{N}_{SP}(\mathbf{0}, \sigma_{\text{spp}}^2 \mathbf{I}_P \otimes \mathbf{\Sigma}_{\text{spp}})$	Setting (a)
(2) spatial	$\mathcal{N}_{SP}(\mathbf{0}, \sigma_{\text{plot}}^2 \mathbf{\Sigma}_{\text{plot}} \otimes \mathbf{I}_S)$	Setting (b)
(3) spatio-phylogenetic	$\mathcal{N}_{SP}(\mathbf{0}, \sigma_{\text{plot}}^2 \sigma_{\text{spp}}^2 \mathbf{\Sigma}_{\text{plot}} \otimes \mathbf{\Sigma}_{\text{spp}})$	Setting (c)
(4) independent	$\mathcal{N}_{SP}(\mathbf{0}, \sigma_{\text{spp}}^2 \mathbf{I}_P \otimes \mathbf{I}_S)$	

Table 1: Distributions used to simulate the random effects for the different correlation structures in each setting (a) through (c).

	RESPONSE TYPE	SETTING (A) PHYLOGENETIC	SETTING (B) SPATIAL	SETTING (C) SPATIO- PHYLOGENETIC
CCM	Gaussian	910	802	900
	binary	-	-	1
	Poisson	1	-	2
ESM	Gaussian	1,308	1,246	1,402
	binary	4	2	10
	Poisson	114	263	254

Table 2: The number of discarded runs in the simulation study (from a total of 4,900 runs per response type, setting and model).

RESPONSE TYPE	β_{spp}	σ_{spp}^2	σ_{plot}^2
Gaussian	$[-2 \quad -2/3 \quad 2/3 \quad 2]^\top$	2	2
binary	$[-1 \quad -1/3 \quad 1/3 \quad 1]^\top$	1	1
Poisson	$[-1 \quad -1/3 \quad 1/3 \quad 1]^\top$	1	1

Table 3: Fixed effects parameters β_{spp} for each of the four species and phylogenetic and spatial random effect variance parameters, σ_{spp}^2 and σ_{plot}^2 respectively, used for each response type in the simulation study of the PSGLMs.

RANDOM EFFECT COMPONENT	TENDENCY	VARIANCE-COVARIANCE MATRIX
spatial		$\sigma_\gamma^2 \Sigma_{\text{plot}} \otimes \mathbf{I}_S$
phylogenetic	attraction	$\sigma_\gamma^2 \mathbf{I}_P \otimes \Sigma_{\text{spp}}$
spatio-phylogenetic	attraction	$\sigma_\gamma^2 \Sigma_{\text{plot}} \otimes \Sigma_{\text{spp}}$
phylogenetic	repulsion	$\sigma_\gamma^2 \mathbf{I}_P \otimes \Sigma_{\text{spp}}^{-1}$
spatio-phylogenetic	repulsion	$\sigma_\gamma^2 \Sigma_{\text{plot}} \otimes \Sigma_{\text{spp}}^{-1}$
independent		$\sigma_\gamma^2 \mathbf{I}_P \otimes \mathbf{I}_S$

Table 4: Variance-covariance matrices for each random effect component used to model the bird abundance data. For a *repulsion* tendency, where closely related species are less likely to co-occur, contrary to the attraction tendency, the inverse of the phylogenetic correlation matrix was used.

ENVIRONMENTAL FACTOR	CHOSEN CORRELATION STRUCTURE	TENDENCY	$\hat{\sigma}_\gamma^2$
elevation	spatial		< 0.001
broadleaved trees	spatio-phylogenetic	attraction	0.007
rejuvenation	spatio-phylogenetic	attraction	0.004
DBH	spatial		< 0.001
tree holes	spatio-phylogenetic	repulsion	0.030

Table 5: For each of five environmental factors, a separate ESM was fitted. With the spatial correlation structure chosen for DBH and elevation, we concluded that no phylogenetic signal is present in the response of the species to the factors. For the other factors, a spatio-phylogenetic signal was detected with a repulsion tendency for tree holes and an attraction tendency for broadleaved trees and rejuvenation. $\hat{\sigma}_\gamma^2$ is the random effect variance estimate for the correlation structure chosen using a predictive cross-validation.

783 List of Figures

784	1	Simulation study for the CCM with Gaussian responses: the	
785		distribution of the variance parameter estimates for each ran-	
786		dom effect component (rows) and setting (columns) subject to	
787		the phylogenetic correlation ρ between species 1 and 2. The	
788		true parameter to be estimated was $\hat{\sigma}_\gamma^2 = 2$ in settings (a) and	
789		(b) and $\hat{\sigma}_\gamma^2 = 2 \times 2 = 4$ in setting (c). Outliers are not drawn.	52
790	2	Selection rates for the CCM. The rate at which the correct ran-	
791		dom effect component (red line) was chosen by the CCM for	
792		each response type and setting (a) through (c) subject to the	
793		phylogenetic correlation ρ between species 1 and 2. The <i>cor-</i>	
794		<i>rect component</i> is the component that was used to simulate the	
795		data in each setting (e.g. phylogenetic component in setting	
796		(a)). The blue line depicts the rate at which any of the random	
797		effect components (phylogenetic or spatio-phylogenetic) was	
798		chosen by the CCM and thus detects a (spatio-)phylogenetic	
799		signal. In setting (b), this rate thus acts as a false-positive	
800		rate and should be low.	53
801	3	Selection rates for the ESM. The rate at which the correct ran-	
802		dom effect component (red line) was chosen by the ESM for	
803		each response type and setting (a) through (c) subject to the	
804		phylogenetic correlation ρ between species 1 and 2. The <i>cor-</i>	
805		<i>rect component</i> is the component that was used to simulate the	
806		data in each setting (e.g. phylogenetic component in setting	
807		(a)). The blue line depicts the rate at which any of the random	
808		effect components (phylogenetic or spatio-phylogenetic) was	
809		chosen by the CCM and thus detects a (spatio-)phylogenetic	
810		signal. In setting (b), this rate thus acts as a false-positive	
811		rate and should be low.	54

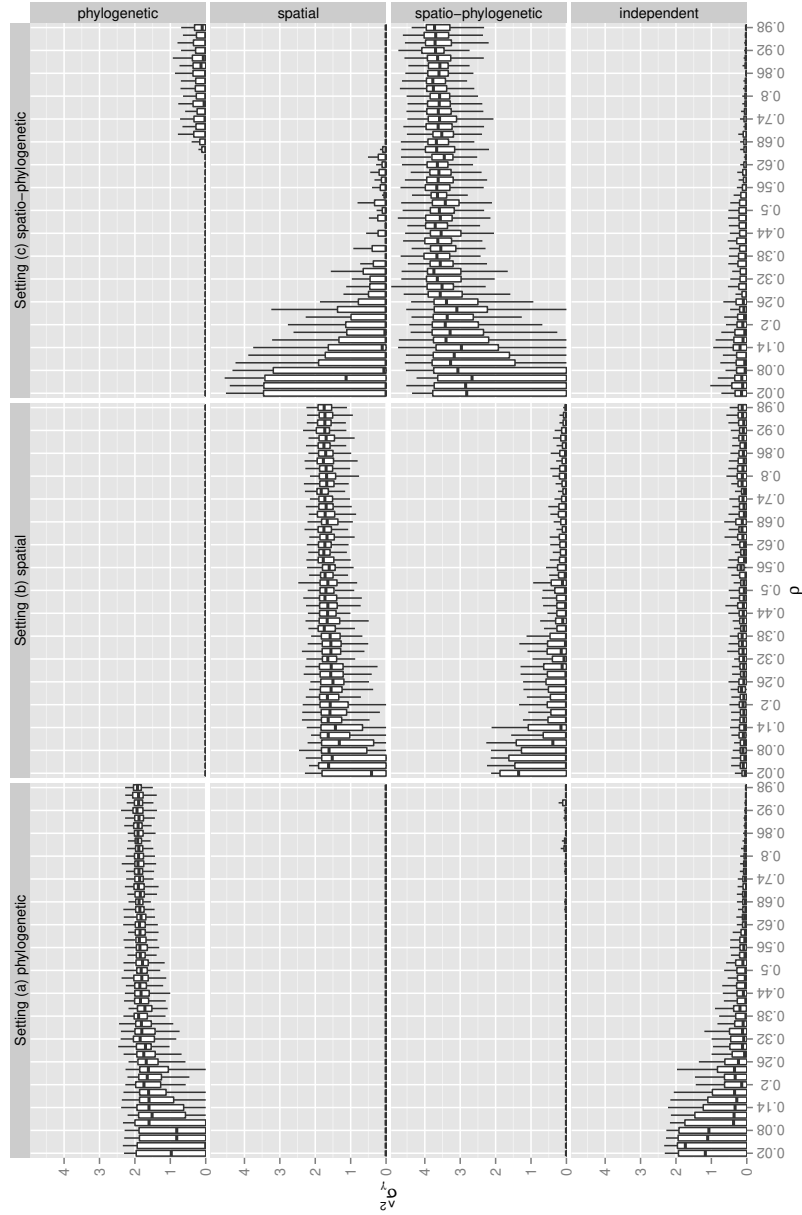


Figure 1: Simulation study for the CCM with Gaussian responses: the distribution of the variance parameter estimates for each random effect component (rows) and setting (columns) subject to the phylogenetic correlation ρ between species 1 and 2. The true parameter to be estimated was $\hat{\sigma}_\gamma^2 = 2$ in settings (a) and (b) and $\hat{\sigma}_\gamma^2 = 2 \times 2 = 4$ in setting (c). Outliers are not drawn.

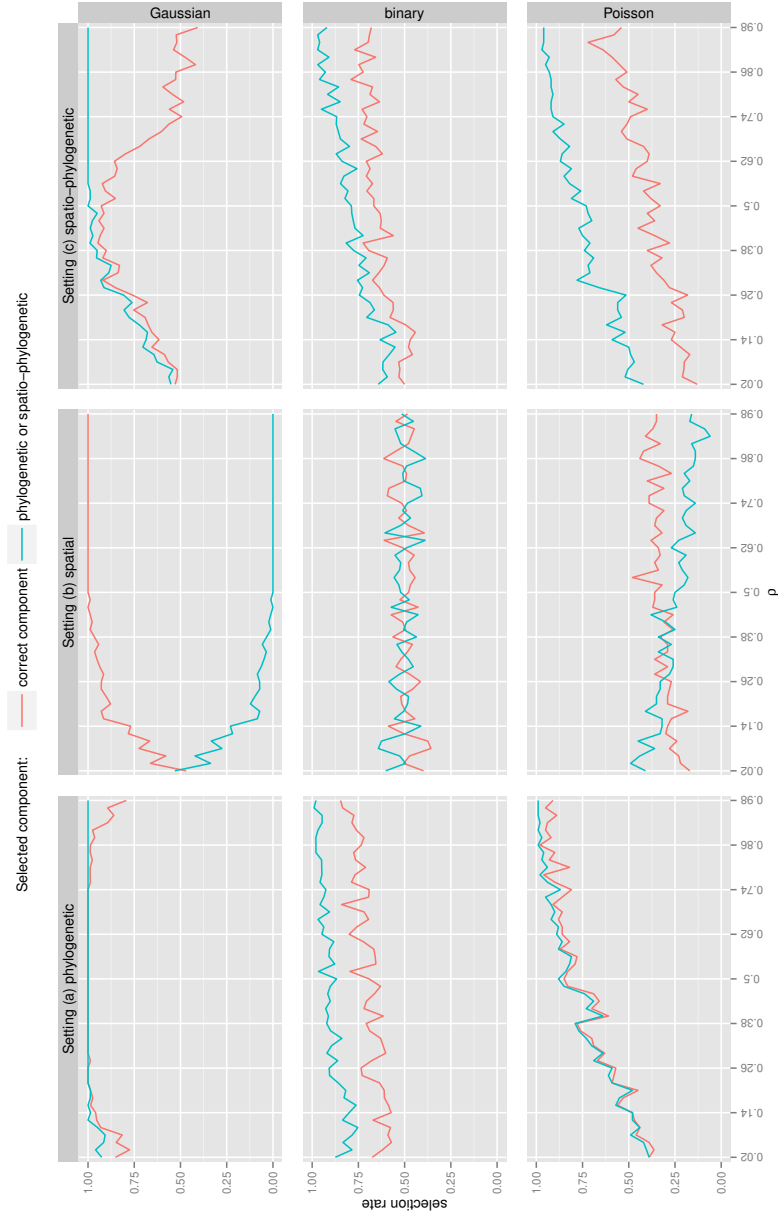


Figure 2: Selection rates for the CCM. The rate at which the correct random effect component (red line) was chosen by the CCM for each response type and setting (a) through (c) subject to the phylogenetic correlation ρ between species 1 and 2. The *correct component* is the component that was used to simulate the data in each setting (e.g. phylogenetic component in setting (a)). The blue line depicts the rate at which any of the random effect components (phylogenetic or spatio-phylogenetic) was chosen by the CCM and thus detects a (spatio-)phylogenetic signal. In setting (b), this rate thus acts as a false-positive rate and should be low.

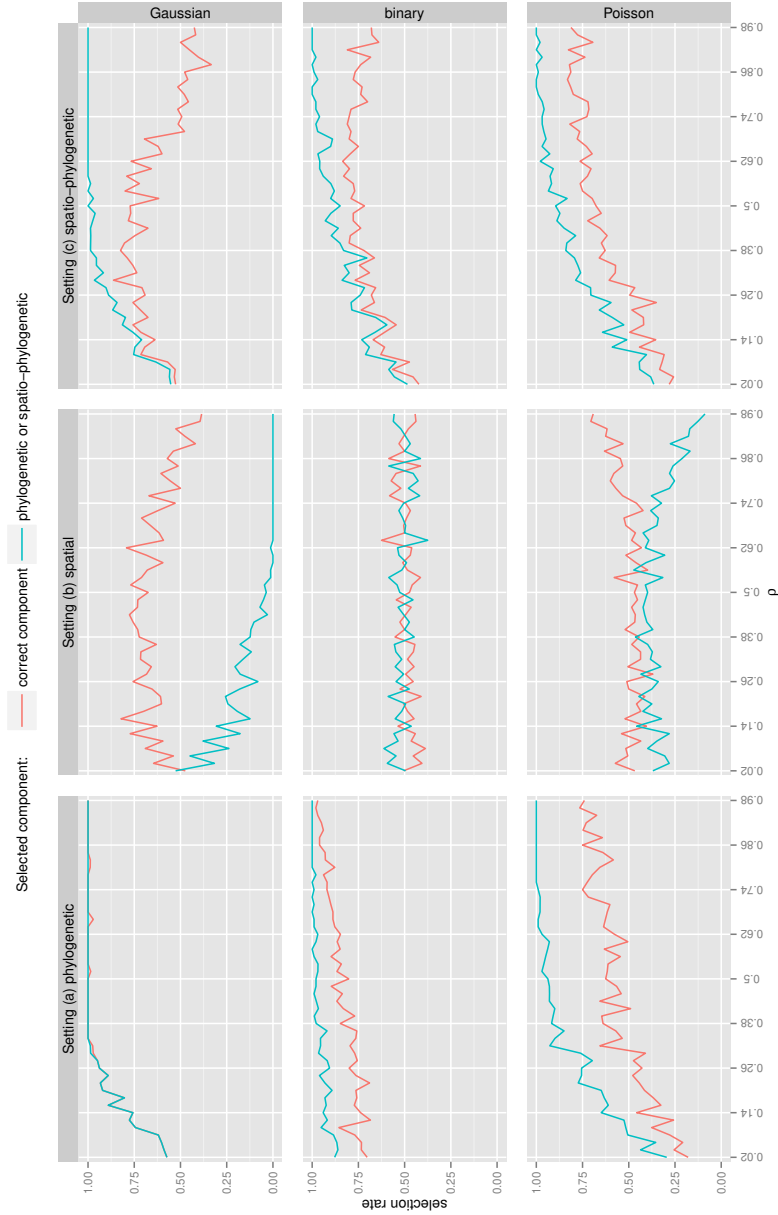


Figure 3: Selection rates for the ESM. The rate at which the correct random effect component (red line) was chosen by the ESM for each response type and setting (a) through (c) subject to the phylogenetic correlation ρ between species 1 and 2. The *correct component* is the component that was used to simulate the data in each setting (e.g. phylogenetic component in setting (a)). The blue line depicts the rate at which any of the random effect components (phylogenetic or spatio-phylogenetic) was chosen by the CCM and thus detects a (spatio-)phylogenetic signal. In setting (b), this rate thus acts as a false-positive rate and should be low.